

BERT

Pre-training of Deep Bidirectional Transformers
For Language Understanding

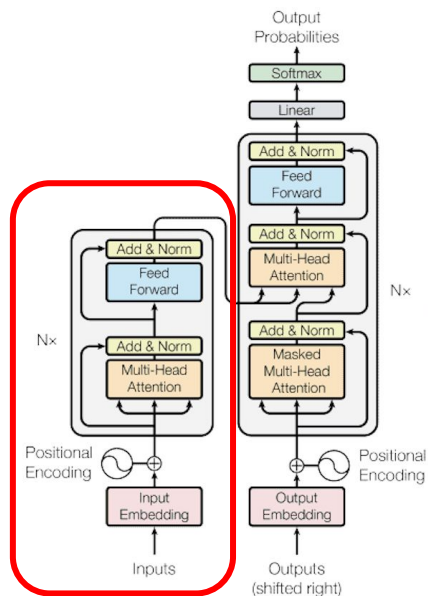


Hunner Markus	01503441
Meier Ronja	12433721
Steinegger Benno	12117772

BERT¹

Bidirectional Encoder Representations from Transformers

» Encoder-Only model

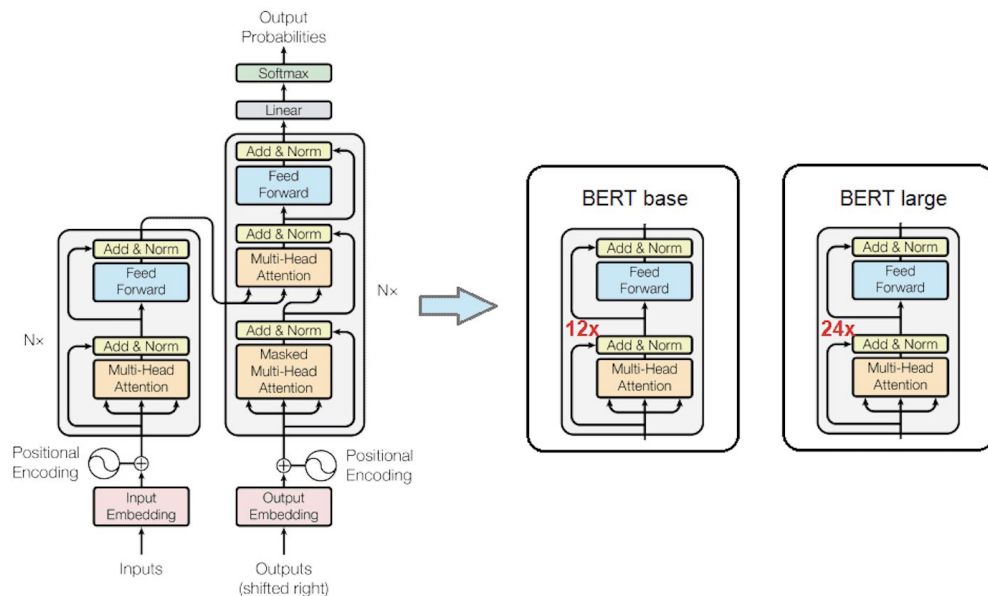


¹[Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding" 2019.](#)

BERT¹

Bidirectional Encoder Representations from Transformers

» Encoder-Only model



¹[Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding" 2019.](#)

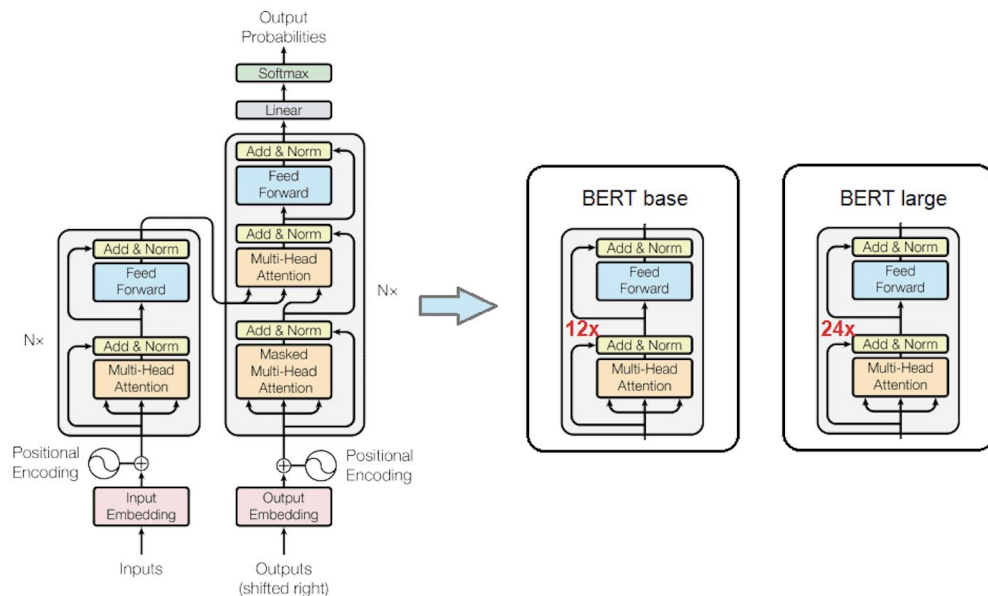
BERT¹

Bidirectional Encoder Representations from Transformers

» Encoder-Only model using

bidirectional attention

- Deep understanding of input text
- Excels at NLP tasks like text classification or question answering



¹[Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding" 2019.](#)

BERT¹

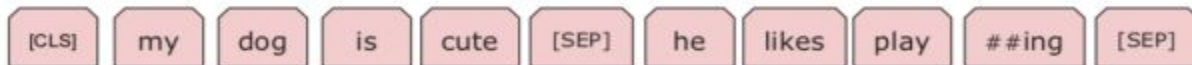
» WordPiece Tokenization

- "unhappiness" → ["un", "##happiness"], "playing" → ["play", "##ing"]
- Small vocabulary → Efficient mapping of text to embeddings & covering of infrequent terms

» Special tokens

- [CLS] Classification token, used as pooling operator to get a single vector per sequence
- [SEP] Used to indicate a second sentence
- [MASK] Used in the masked language model, to predict this word
- [PAD] Added to sequences to ensure all inputs in a batch have the same length

Input



¹[Devlin, Jacob, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding" 2019.](#)

BERT - Model

» Model architecture:

- n layers of stacked transformers
- Every transformer's layer receives as input the output of the previous layer
- Model head depending on problem

» Pre-training & workflow

- **Pre-training** a large model on huge datasets needs a lot of computation
- Download a pre-trained model, attach a head fit for the problem task, and **fine-tune** on own dataset

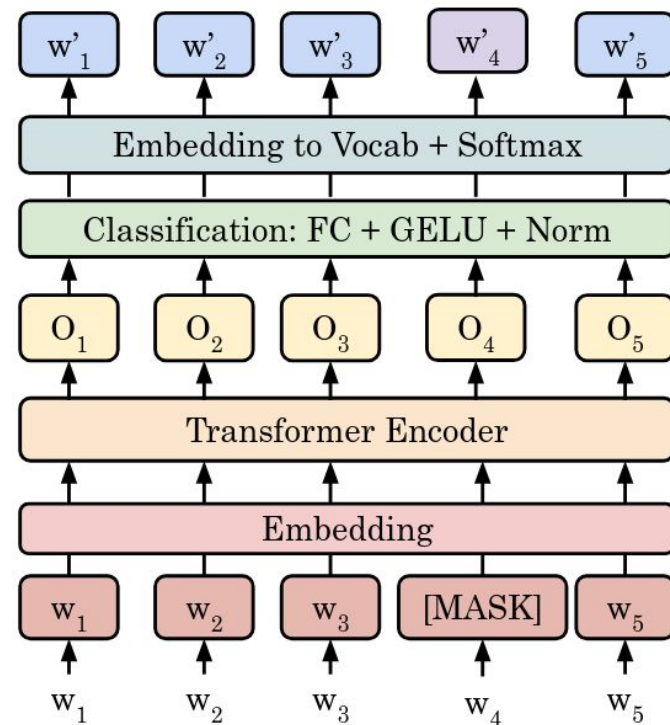


Figure: BERT Model Architecture. Figure 4 (a) in [1]

Pre-Training - Motivation

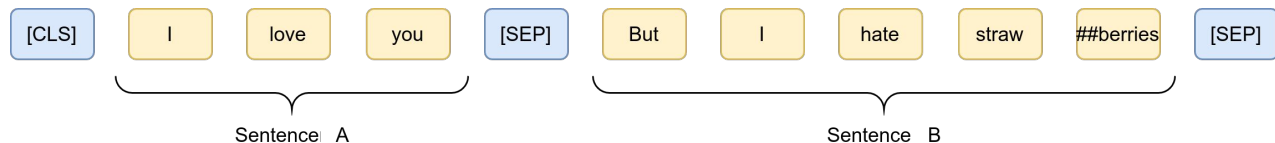
- » For many tasks we do not have much training data
- » Large models need a lot of data to work well

Solution: Train unsupervised on task-agnostic data

- » Train model about meaning of words and patterns in language
- » Train without labels, and instead predict word and sentence positions

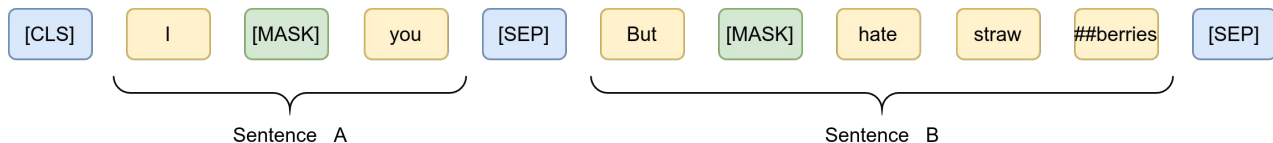
Pre-Training

- » Take any large text corpus ...
- » ... and form sentence pairs
 - 50% of the time B is the actual next sentence that follows A (labeled as *IsNext*), and 50% of the time it is a random sentence from the corpus (labeled as *NotNext*).



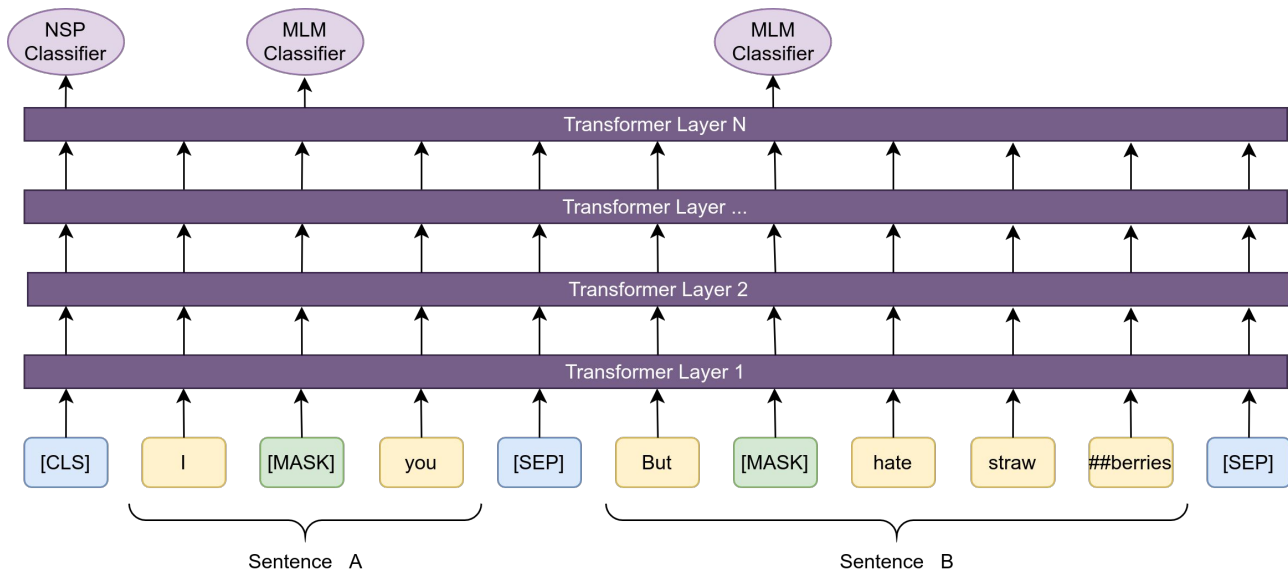
Pre-Training

- » Take any large text corpus ...
- » ... and form sentence pairs
 - 50% of the time B is the actual next sentence that follows A (labeled as *IsNext*), and 50% of the time it is a random sentence from the corpus (labeled as *NotNext*).
- » ... and replace 15% of all tokens by **[MASK]**



Pre-Training

- » ... send it through the model and
- » ... give **[CLS]** token result to Next-Sentence-Prediction (NSP) classifier
- » ... give **[MASK]** token results to Masked-Language-Modeling (MLM) classifier



Pre-Training

- » ... send it through the model and
 - » ... all classifier contribute to the loss. → They do not learn separately!
-
- » Someone with lots of compute and time **pre-trains** a large model ...
 - » ... and we download it and **fine-tune** on our own data.

Fine-Tuning

- » To fine-tune BERT we simply replace the model heads with task-specific ones.
- » For example: Add a Sequence Classification onto the [CLS] token output.
 - We do not need a pooling layer, as we trained the model to encode all necessary information in the [CLS] vector output.
- » Fine-tuning for 2 to 4 epochs is usually sufficient.

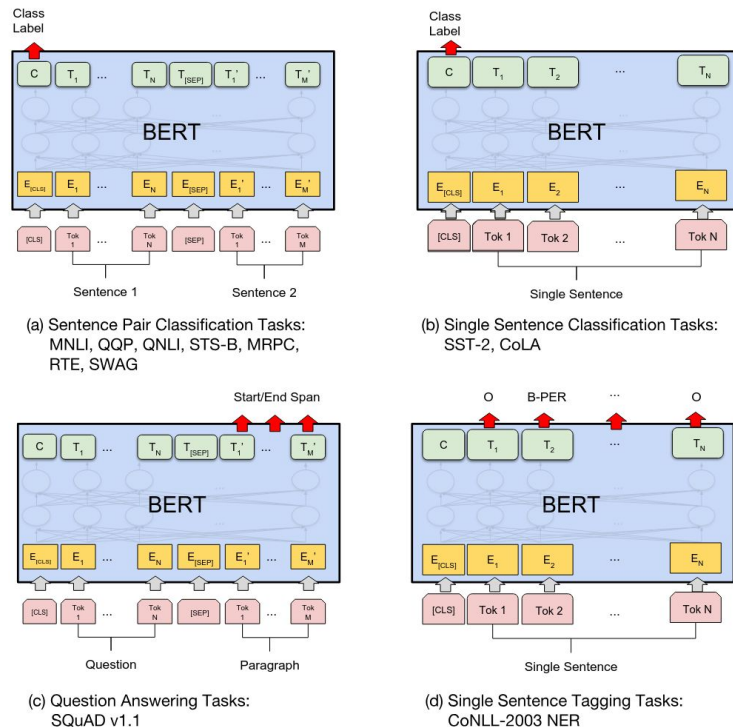


Figure: Illustrations of Fine-tuning BERT on Different Tasks.

Figure 4 in [1]

ModernBERT

Important architectural changes

ModernBERT³ - Structural Changes

» Bias Terms

- Disable bias in all linear layers except the in final decoder linear layer
- Disable in Layer Norms

» Positional Embeddings

- Rotary positional Embeddings (RoPE) instead of absolute positional embeddings

» Activation

- GeGLU based on original BERT's GeLU
- Gated-Linear Units based

» Normalization

- Pre-normalization block & standard layer normalization to stabilize training

[3] B. Warner et al., "Smarter, Better, Faster, Longer: A Modern Bidirectional Encoder for Fast, Memory Efficient, and Long Context Finetuning and Inference," 2024, arXiv. doi: 10.48550/ARXIV.2412.13663.

ModernBERT - Model Design

- » Remove NSP
- » Masked Language Modeling (30% instead of 15%; from MosaicBERT)
- » Tokenizer
 - BPE tokenizer based on OLMo
 - Same special tokens (e.g. [CLS], [SEP], ...)
- » Deeper & Narrow Layer
 - 22 (base), 28 (large) layers
- » Higher Context Length of 8192
- » Learning Rate & Batch Size Scheduler

ModernBERT - Efficiency Improvements

» Attention

- Alternating Attention: alternates between global and a sliding window attention
- Flash Attention 3

» Unpadding

- Remove padding tokens

» Use torch.compile

GLUE Fine-tuning

General Language Understanding Evaluation

GLUE²

General Language Understanding Evaluation

- » A whole suite of **9 sentence classification problems** in English.
 - 6 classification problems with 2 labels (**CoLA**, **MRPC**, **QNLI**, **QQP**, **RTE**, **SST2**)
 - 2 classification problems with 3 labels (**MNLI-m**, **MNLI-mm**)
 - 1 regression problem within values of 1 to 5 (**STSB**)
- » A set of standardized public datasets.
 - 8 training datasets (Problems **MNLI-m**, **MNLI-mm** share a train set)
 - 9 validation datasets
 - 9 tests sets with *censored labels*

²[Wang, Alex, et al. "GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding." 2018.](#)

GLUE²

General Language Understanding Evaluation

» A whole suite of **9 sentence classification problems** in English.

- 6 classification problems with 2 labels (**CoLA, MRPC, QNLI, QQP, RTE, SST2**)
- 2 classification problems with 3 labels (**MNLI-m, MNLI-mm**)
- 1 regression problem within values of 1 to 5 (**STSB**)

» A set of standardized public datasets.

- 8 training datasets (Problems **MNLI-m, MNLI-mm** share a train set)
- 9 validation datasets
- ~~9 tests sets with censored labels~~



We need a labeled test set to provide results ...

GLUE - Results of original Paper

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

“Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks.” [1]

¹[Devlin, Jacob, et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding” 2019.](#)

GLUE - Our Results

System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Turing ULR v6	92.5/92.1	76.4	96.7	97.5	73.3	93.1	94.2	93.6	89.9
Vega v1	92.1/91.9	76.7	96.7	97.9	73.8	93.1	94.5	92.4	89.9
Turing NLR v5	92.6/92.4	76.4	97.9	97.6	72.6	93.3	93.8	94.1	90.0
Models compared to in original paper:									
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
Results of original BERT paper:									
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1
Our reproduction:									
BERT _{BASE}	86.4/84.2	87.7	89.4	94.6	54.8	84.9	85.6	64.7	81.4
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2
Our comparison with ModernBERT:									
ModernBERT _{BASE}	88.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	76.3	84.7
ModernBERT _{LARGE}	89.2/87.6	88.7	91.5	94.2	60.2	90.6	90.5	83.1	86.2

GLUE - Our Results

System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Turing ULB v0	92.5/92.1	76.4	96.7	97.5	73.3	93.1	94.2	93.6	90.9
Vega v1	92.1/91.9	76.7	96.7	97.5	73.6	93.1	94.5	92.1	90.9
Turing NLB v0	92.6/92.4	76.4	97.8	97.6	73.6	93.2	93.9	94.1	90.9
Models compared to in original paper:									
Pre-OpenAI SOTA	90.6/90.1	66.1	92.3	93.2	35.0	81.0	90.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	92.1/91.4	70.3	97.4	91.3	45.4	90.0	92.3	56.0	75.1
Results of original BERT paper:									
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1
Our reproduction:									
BERT _{BASE}	86.4/84.2	87.7	89.4	94.6	54.8	84.9	85.6	64.7	81.4
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2
Our comparison with ModernBERT:									
ModernBERT _{BASE}	86.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	76.3	84.7
ModernBERT _{LARGE}	89.2/87.6	88.7	91.5	94.2	60.2	90.6	90.5	83.1	86.2

GLUE - Our Results

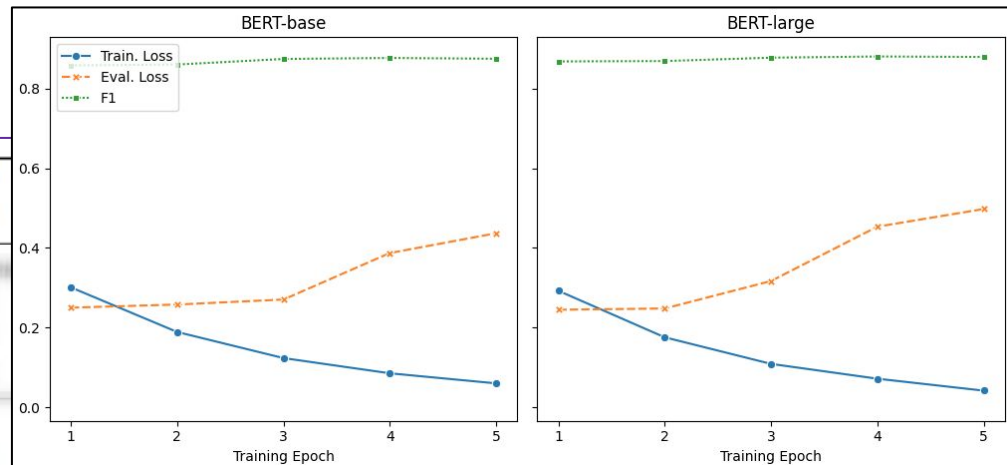
System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Training GLUE v0	92.5/92.1	76.4	96.7	97.3	73.3	95.1	94.2	93.6	89.9
Yogi v1	92.1/91.9	76.7	96.7	97.4	73.4	95.1	94.5	93.1	89.9
Overall, we achieve comparable prediction performance.									
Pre-OpenAI SOFA	90.6/90.1	66.1	92.3	93.2	61.1	91.1	90.7	91.7	74.9
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.6	90.1	36.0	71.1	76.8	71.9	71.9
OpenAI GPT	92.1/91.4	70.3	97.4	98.3	45.4	90.0	92.1	91.0	75.1
Results of original BERT paper:									
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1
Our reproduction:									
BERT _{BASE}	86.4/84.2	87.7	89.4	94.6	54.8	84.9	85.6	64.7	81.4
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2



Not a valid comparison: We do not evaluate on the same test dataset!

GLUE - Our Results

System	QQP dataset size	QQP
Original Train:	363k	363k
Original Test:	390k	
Our Train:	291k	
Our Test:	72k	



Models compared to in original paper:										
Pro-OpenAI GPT-3	80.6/80.1	80.1								
BLSTM+ELMo+Attn	76.4/76.1	64.8	79.6	90.4	96.0	73.3	84.9	96.8	71.0	
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	95.4	80.0	82.3	96.0	75.1	
Results of original BERT paper:										
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6	
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1	
Our reproduction:										
BERT _{BASE}	86.4/84.2	87.7	89.4	94.6	54.8	84.9	85.6	64.7	81.4	
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2	




Not a valid comparison: We do not evaluate on the same test dataset!



GLUE - Our Results

System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Turing ULB v0	92.5/92.1	76.4	96.7	97.3	73.3	93.1	94.2	93.6	90.9
Vega v1	92.1/91.9	76.7	96.7	97.3	73.4	93.1	94.5	92.1	90.9
Turing NLB v1	92.6/92.4	76.1	97.0	97.6	72.6	93.2	93.9	94.1	90.9
Models compared to in original paper:									
Pre-OpenAI SOTA	90.6/90.1	66.1	92.3	93.2	35.0	81.0	90.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	76.8	71.0
OpenAI GPT	92.1/91.4	70.3	97.4	98.3	45.4	90.0	92.3	76.0	75.1
Results of original BERT paper:									
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1
Our reproduction:									
BERT _{BASE}	86.4/84.2	87.7	89.4	94.6	54.8	84.9	85.6	64.7	81.4
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2
Our comparison with ModernBERT:									
ModernBERT _{BASE}	88.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	76.3	84.7
ModernBERT _{LARGE}	89.2/87.6	88.7	91.5	94.2	60.2	90.6	90.5	83.1	86.2

GLUE - Our Results

System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Training ULB v0	92.5/92.1	76.4	96.7	97.3	73.3	93.1	94.2	93.6	90.9
Vega v1	92.1/91.9	76.7	96.7	97.3	73.4	93.1	94.5	92.1	90.9
Training NLI v1	92.4/92.4	76.1	97.3	97.1	73.1	93.3	93.9	94.1	90.9
Models compared to in original paper:									
 ModernBERT_{BASE} achieves prediction performance comparable to BERT_{LARGE}									
OpenAI GPT-1	92.1/91.4	78.3	97.4	95.3	49.4	90.0	92.3	96.0	79.1
Results of original BERT paper:									
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1
Our reproduction:									
BERT _{BASE}	86.4/84.2	87.7	89.4	94.6	54.8	84.9	85.6	64.7	81.4
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2
Our comparison with ModernBERT:									
ModernBERT _{BASE}	88.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	76.3	84.7
ModernBERT _{LARGE}	89.2/87.6	88.7	91.5	94.2	60.2	90.6	90.5	83.1	86.2

GLUE - Our Results

System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Training ULB v0	82.5/82.1	76.4	86.7	87.3	73.3	85.1	84.2	83.6	80.9
Vega v1	82.1/81.9	76.7	86.7	87.8	73.4	85.1	84.5	82.1	80.9
Training NLB v0	82.4/82.4	76.1	87.8	87.6	73.6	85.2	83.9	84.1	80.9
Models compared to in original paper:									
 ModernBERT_{BASE} achieves prediction performance comparable to BERT_{LARGE}									
OpenAI GPT-2	82.1/81.4	76.2	87.4	85.3	49.4	80.0	82.3	86.0	75.1
 ModernBERT_{BASE} achieves significantly better runtime performance than BERT_{BASE}									
Our reproduction:									
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2
Our comparison with ModernBERT:									
ModernBERT _{BASE}	88.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	76.3	84.7
Speedup in training time									
BERT _{BASE} / ModernBERT _{BASE}	1.38	1.24	1.38	0.99	0.76	1.15	1.18	1.36	1.18

GLUE - Our Results

System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Turing ULR v6	92.5/92.1	76.4	96.7	97.5	73.3	93.1	94.2	93.6	89.9
Vega v1	92.1/91.9	76.7	96.7	97.9	73.8	93.1	94.5	92.4	89.9
Turing NLR v5	92.6/92.4	76.4	97.9	97.6	72.6	93.3	93.8	94.1	90.0
Models compared to in original paper:									
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
Results of original BERT paper:									
BERT _{base}	84.4/83.4	71.2	90.3	90.3	52.1	85.8	88.9	88.4	79.6
BERT _{Large}	86.7/85.9	72.1	92.7	91.9	60.5	86.5	89.3	79.1	82.1
Our reproduction:									
BERT _{base}	86.4/84.2	87.7	90.4	94.6	54.8	84.9	85.6	84.7	81.4
BERT _{Large}	88.2/85.3	88.2	90.3	94.7	58.8	87.6	88.1	87.3	83.2
Our comparison with ModernBERT:									
ModernBERT _{base}	88.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	78.3	84.7
ModernBERT _{Large}	89.2/87.6	88.7	91.5	94.2	60.2	90.6	90.5	83.1	86.2

GLUE - Our Results

System	MNLI (m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Avg. -
Current Top 3 Models according to GLUE leaderboard:									
Turing ULR v6	92.5/92.1	76.4	96.7	97.5	73.3	93.1	94.2	93.6	89.9
Vega v1	92.1/91.9	76.7	96.7	97.9	73.8	93.1	94.5	92.4	89.9
Turing NLR v5	92.6/92.4	76.4	97.9	97.6	72.6	93.3	93.8	94.1	90.0
Models compared to in original paper:									
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
Results of original BERT									
BERT _{base}	88.5/88.1	70.2	89.2	92.1	39.1	85.4	89.1	67.2	78.4
BERT _{large}	90.4/90.1	71.2	90.2	93.3	40.2	86.5	90.3	68.7	80.0
Our comparison with ModernBERT:									
ModernBERT _{base}	90.4/90.0	70.3	90.6	94.1	39.9	89.0	89.7	76.3	84.7
ModernBERT _{large}	90.2/87.6	88.7	91.5	94.2	40.2	90.6	90.5	83.1	86.2



Many of the task in the GLUE benchmark can not be considered particularly “hard” problems anymore.

SQuAD Fine-tuning

Stanford Question Answering Dataset

SQuAD dataset

- Reading comprehension dataset
- Questions posed by crowdworkers on a set of Wikipedia articles
- Training dataset: 87599 | Validation dataset: 10570
- Focus on SQuAD 1.1 task
 - Extract the correct answer from a context to a given question
- SQuAD 2.0
 - Extends 1.1 by considering the possibility that no short answer exists in the provided paragraph
→ more realistic

SQuAD 1.1 results of the BERT paper

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
Published				
BiDAF+ELMo (Single)	-	85.6	-	85.8
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

SQuAD 1.1 results of the BERT paper

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
Published				
BiDAF+ELMo (Single)	-	85.6	-	85.8
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.0
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

Paper does not specify which 7 pre-training checkpoints and which fine-tuning seeds are used

TriviaQA problem with Training loss

SQuAD 1.1 results of the BERT paper

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
Published				
BiDAF+ELMo (Single)	-	85.6	-	85.8
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT_{LARGE} (Ensemble)	85.8	91.8		
BERT_{LARGE} (Sgl. + TriviaQA)	84.2	91.1	85.1	91.8
BERT_{LARGE} (Ens. + TriviaQA)	86.2	92.2	87.4	93.2

SQuAD 1.1 our results

Model	EM		F1	
	Paper	Our	Paper	Our
BERT-Base	80.8	79.9	88.5	87.7
BERT-Large	84.1	83.3	90.9	90.5

SWAG Fine-tuning

A Large-Scale Adversarial Dataset for
Grounded Commonsense Inference

SWAG dataset

- » **S**ituations **W**ith **A**dversarial **G**enerations.
- » designed for grounded common sense inference
- » Predict next plausible sequence/situation
 - Generated from video captions
- » Multiple-choice format (1 out of 4 is correct)
 - AutoModelForMultipleChoice, own class for ModernBERT and GPT-2
 - Dropout, Linear Classifier (hidden dimension \rightarrow 4) and Soft max
- » Adversarial Filtering

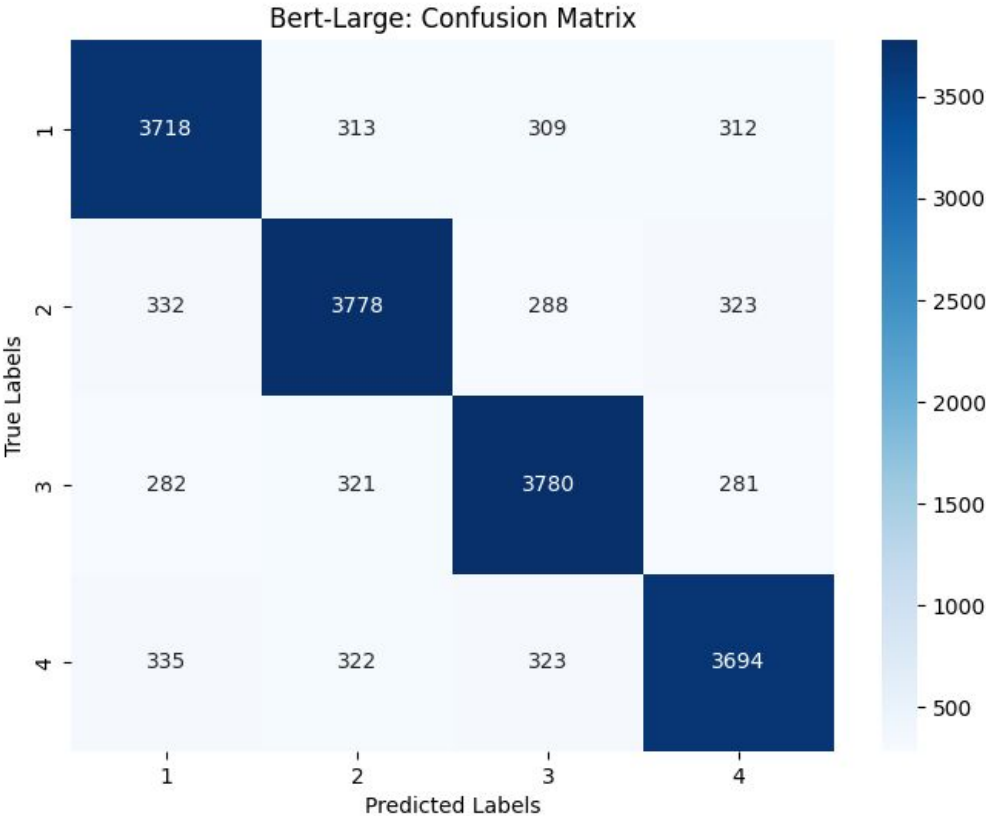
SWAG results

Model	Validation		Test		Parameters
	Paper	Our	Paper	Our	
BERT-Base	81.6	76.6	-	76.1	110M
BERT-Large	86.6	79.3	86.3	80.8	340M
OpenAI GPT	-	66.8	78.0	66.2	117M
ModernBERT-base	-	77.6	-	77.8	149M
ModernBERT-base _{xavier}	-	77.7	-	78.0	149M
ModernBERT-large	-	82.0	-	82.1	395M
DeBERTa-v3-base	-	83.6	-	83.4	100M

SWAG

Model

BERT-Bas
BERT-Lar
OpenAI G
ModernBE
ModernBE
ModernBE
DeBERTa-




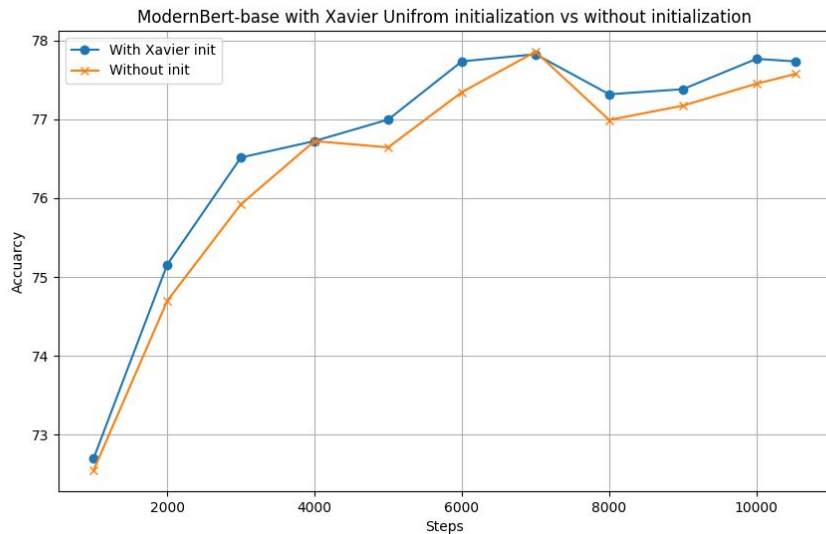
aining
plit,
e not

ters

110M
340M
117M
149M
149M
395M
100M

SWAG results

Model	Validation		Test		Parameters
	Paper	Our	Paper	Our	
BERT-Base	81.6	76.6	-	 We used a different and smaller version of GPT	100M
BERT-Large	86.6	79.3	86.3		340M
OpenAI GPT	-	66.8	78.0	66.2	117M
ModernBERT-base	-	77.6	-	77.8	149M
ModernBERT-base _{xavier}	-	77.7	-	78.0	149M
ModernBERT-large	-	82.0	-	82.1	395M
DeBERTa-v3-base	-	83.6	-	83.4	100M



Test		Parameters
Paper	Our	
-	76.1	110M
86.3	80.8	340M
78.0	66.2	117M
-	77.6	149M
-	77.7	149M
-	82.1	395M
-	3.4	100M

ModernBERT-base

ModernBERT-base_{xavier}


ModernBERT-large

DeBERTa-v3-base



We get slightly better results with Xavier-weight-initialization

SWAG results

Model	Validation		Test		Parameters
	Paper	Our	Paper	Our	
BERT-Base	81.6	76.6	-	76.1	110M
BERT-Large	86.6	79.3	86.3	80.8	340M
OpenAI GPT	-	66.8	78.0	66.2	117M
 Current SOTA: DeBERTa-Large with 94% : Uses disentangled attention and an enhanced mask decoder				77.8	149M
				78.0	149M
ModernBERT-Large	-	82.0	-	82.1	395M
DeBERTa-v3-base	-	83.6	-	83.4	100M

SNLI Fine-tuning

Stanford Natural Language Inference

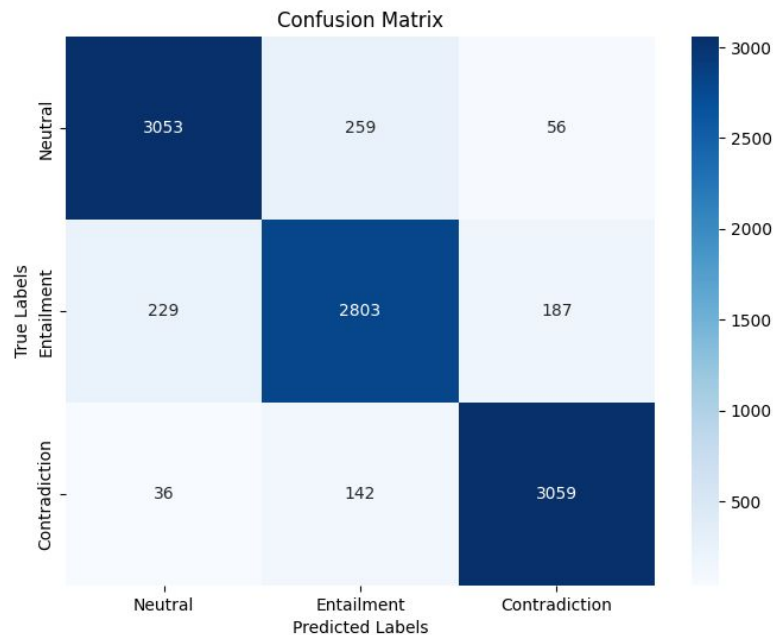
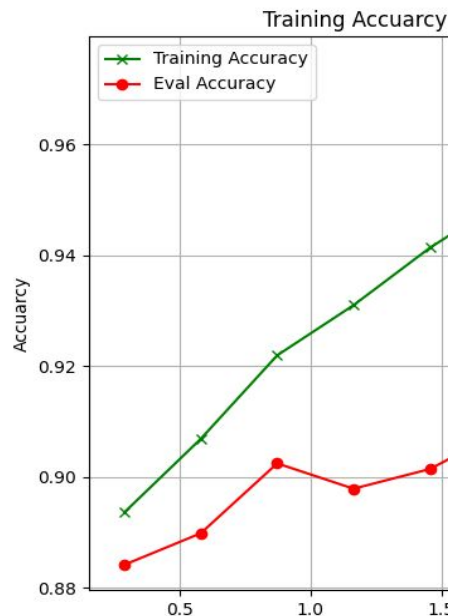
SNLI dataset

- » **S**tanford **N**atural **L**anguage **I**nference
- » 570k human written English sentence pairs
- » Text + Hypothesis → Judgement (Contradiction, Neutral, Entailment)
- » Custom Class *NLIClassifier*
 - Adds a linear classification layer after BERT
 - Would have been also possible with *AutoModelForSequenceClassification*

SNLI Results

» Very high accuracy for BERT_{BASE}: 90.4%

⚠ Due to the size (570k) and high performance of BERT_{BASE}, we didn't finetune BERT_{LARGE}



Evaluation Loss



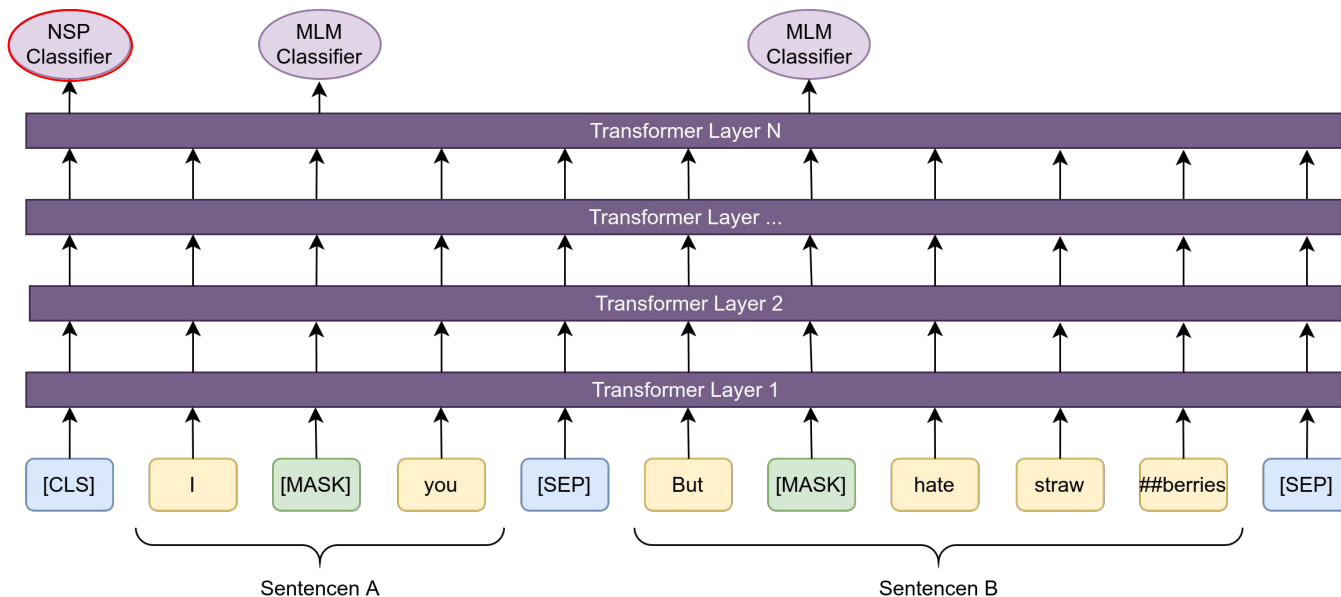
Thank you for your attention.

Implementation of experiments:

<https://github.com/markhun/2024W-DLNLBP-BERT>

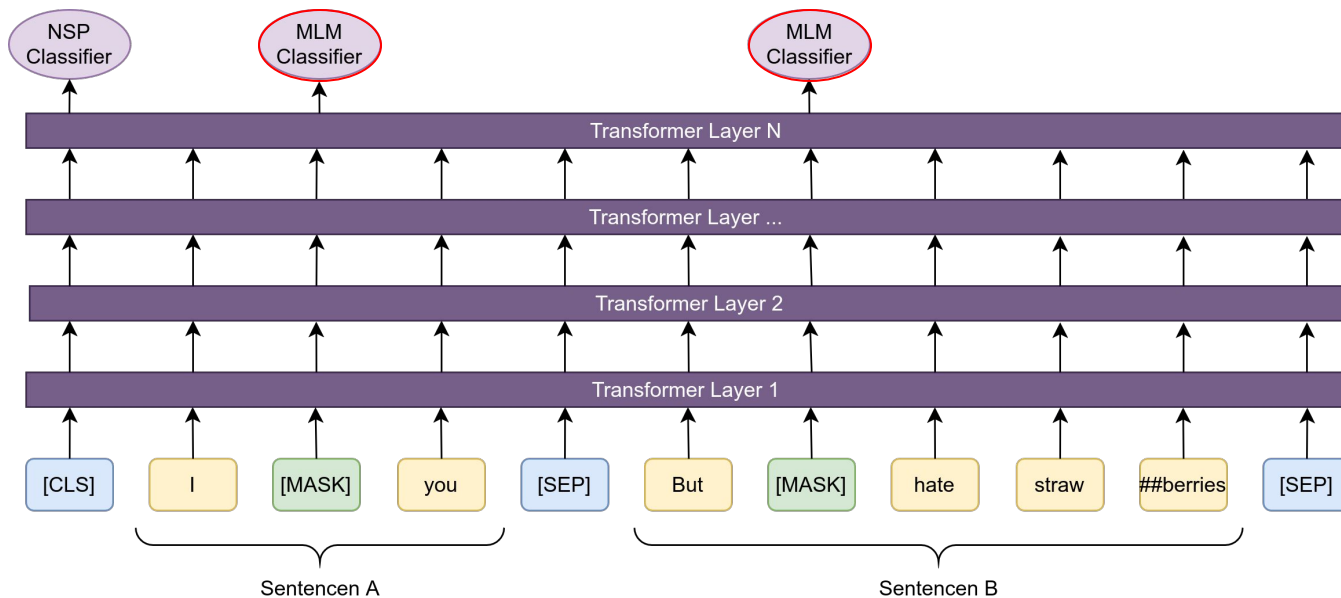
Pre-Training

- » ... send it through the model and
- » ... give [CLS] token result to Next-Sentence-Prediction (NSP) classifier



Pre-Training

- » ... send it through the model and
- » ... give [MASK] token results to Masked-Language-Modeling (MLM) classifier



GLUE Fine-tuning (extended)

General Language Understanding Evaluation

GLUE²

General Language Understanding Evaluation

» A whole suite of 11 sentence classification problems in English.

- 1 Diagnostic Set (**AX**)
- 7 classification problems with 2 labels (**CoLA**, **MRPC**, **QNLI**, **QQP**, **RTE**, **SST2**, **WNLI**)
- 2 classification problems with 3 labels (**MNLI-m**, **MNLI-mm**)
- 1 regression problem within values of 0 to 5 (**STSB**)

» A set of standardized public datasets.

- 9 training datasets (Problems **MNLI-m**, **MNLI-mm** and **AX** share the same training dataset)
- 10 validation datasets (No validation dataset provided for **AX**)
- 11 tests sets with censored labels

²[Wang, Alex. et al. "GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding." 2018.](#)

GLUE²

General Language Understanding Evaluation

» A whole suite of 11 sentence classification problems in English.

- ~~1 Diagnostic Set (AX)~~
- 6 ~~7~~ classification problems with 2 labels (CoLA, MRPC, QNLI, QQP, RTE, SST2, ~~WNLI~~)
- 2 classification problems with 3 labels (MNLI-m, MNLI-mm)
- 1 regression problem within values of 0 to 5 (STSB)

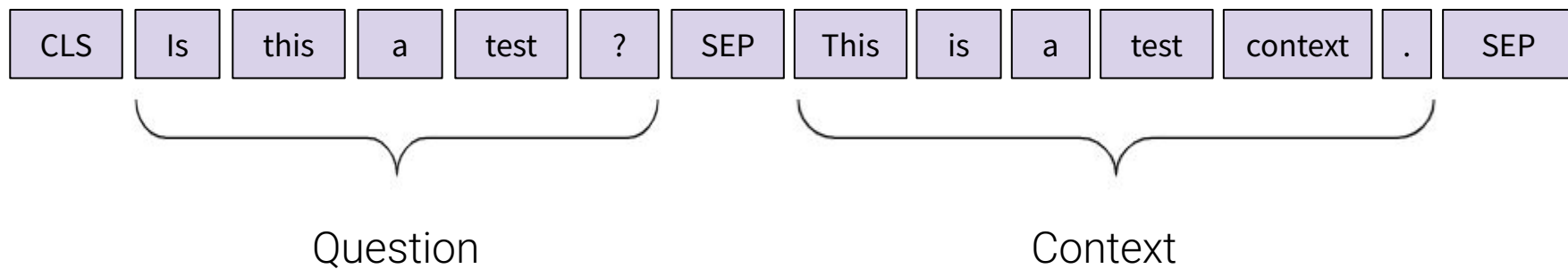
⚠ Considered problematic as it includes adversarial items

» A set of standardized public datasets.

- 9 training datasets (Problems MNLI-m, MNLI-mm and AX share the same training dataset)
- 10 validation datasets (No validation dataset provided for AX)
- 11 tests sets with censored labels

²[Wang, Alex, et al. "GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding." 2018.](#)

Input sequence for SQuAD 1.1



Example of the training dataset

```
{  
  "answers": {  
    "answer_start": [1],  
    "text": ["This is a test text"]  
  },  
  "context": "This is a test context.",  
  "id": "1",  
  "question": "Is this a test?",  
  "title": "train test"  
}
```

- Validation dataset allows several possible answers for each sample